



TITLE:

# Constrained Markov Decision Processes : The Average Case(Mathematical Structure of Optimization Theory)

AUTHOR(S):

Kurano, Masami; Huang, Youqiang

---

CITATION:

Kurano, Masami ...[et al]. Constrained Markov Decision Processes : The Average Case(Mathematical Structure of Optimization Theory). 数理解析研究所講究録 1994, 864: 1-8

ISSUE DATE:

1994-04

URL:

<http://hdl.handle.net/2433/83914>

RIGHT:

# Constrained Markov Decision Processes: The Average Case

蔵野 正美 (Masami Kurano) 黄 佑強 (Youqiang Huang)  
千葉大学教育学部 (Chiba University)

## 1. Preface

A Markov decision process (MDP) with multiple constraints is applicable in many fields. Recently, many researchers pay attention to it ([2],[10] and their references). The methods of analysis about solving it are by using Lagrange multiplier theory ([4],[5]) or by changing it to linear programming (LP) to apply LP theory ([11],[12],[13]).

In this report we consider the average reward criterion of MDP with sample path constraints by the above two ideas under state and action sets being compact. we prove the existence theorem of a constrained optimal pair. Also, by introducing the concept of state-wise mixed policy, we give a characterization of it.

## 2 . Formulation

In this report mentioned Borel sets are Borel subsets of a complete separable metric space. For a Borel set  $X$ ,  $B_X$  denotes the Borel subsets of  $X$ .  $C(X)$  denotes the set of all bounded continuous functions on  $X$ . A Markov decision process with multiple constraints is a controlled dynamic system defined by following objects:  $S$ ,  $\{A(x), x \in S\}$ ,  $Q$ ,  $r$ ,  $c_i$  ( $i = 1, \dots, k$ ), where  $S$  is any Borel set representing the state space of some system and for each  $x \in S$ , the admissible action space  $A(x)$  is a non-empty subset of some Borel set  $A$  such that  $\{(x, a) : x \in S, a \in A(x)\}$  is an element of  $B_S \times B_A$ , the immediate reward function  $r$  is a real-valued Borel measurable function on  $S \times A$ , the immediate cost functions  $c_i$  ( $i=1, \dots, k$ ) are real-valued cost functions on  $S \times A$ ,  $Q(\cdot|x, a)$  is the law of motion, which is taken to be stochastic kernel on  $B_S \times S \times A$ ; i.e, for each  $(x, a) \in S \times A$ ,  $Q(\cdot|x, a)$  is a probability measure on  $B_S$ ; and for each  $D \in B_S$ ,  $Q(D|\cdot)$  is a Borel measurable function on  $S \times A$ .

Throughout this report, the following assumptions will be remain operative:

- (i).  $S$  and  $\{A(x), x \in S\}$  are compacts;
- (ii).  $r$  is non-negative bounded continuous;
- (iii).  $c_i$  is non-negative bounded continuous ( $i = 1, \dots, k$ );
- (iv). whenever  $x_n \rightarrow x, a_n \rightarrow a$ ,  $Q(\cdot|x_n, a_n)$  converges weekly to  $Q(\cdot|x, a)$ .

The sample space is the product space  $\Omega = (S \times A)^\infty$  such that the projections  $X_t, \Delta_t$  on the  $t$ th factors  $S, A$  describe the state and action of the  $t$ th time of the process ( $t \geq 0$ ).

A policy is a sequence  $\pi = (\pi_0, \pi_1, \dots)$  such that, for each  $t \geq 0$ ,  $\pi_t$  is a stochastic kernel on  $B_A \times S \times (A \times S)^t$  with  $\pi_t(A(x_t)|x_0, a_0, \dots, a_{t-1}, x_t) = 1$  for all  $(x_0, a_0, \dots, a_{t-1}, x_t) \in S \times (A \times S)^t$ .

Let  $\Pi$  denote the class of policies.

$T(A | S)$  is the set of all stochastic kernels  $\Phi$  on  $B_A \times S$  with  $\Phi(A(x)|x) = 1$  for all  $x \in S$ .

A policy  $\pi = (\pi_0, \pi_1, \dots)$  is a randomize stationary policy if there is a  $\Phi \in T(A | S)$  such that  $\pi_t(\cdot | x_0, a_0, \dots, x_t) = \Phi(\cdot | x_t)$  for all  $(x_0, a_0, \dots, x_t) \in S \times (A \times S)^t$  and  $t \geq 0$ . Let  $\Phi^\infty$  denote the corresponding policy.

For any  $D \in B_S$ ,  $B(D \rightarrow A)$  denotes the set of all Borel measurable functions  $u: D \rightarrow A$  with  $u(x) \in A(x)$  for all  $x \in D$ .

A randomize stationary policy  $\Phi^\infty$  is called stationary if there is an  $f \in B(S \rightarrow A)$  such that  $\Phi(f(x)|x) = 1$  for all  $x \in S$ . Such a policy will be written by  $f^\infty$ .

$\Pi_{RS}$  and  $\Pi_S$  are respectively the sets of all randomize stationary and stationary policies.

Let  $H_t = (X_0, \Delta_0, \dots, \Delta_t, X_t)$ . It is assumed that, for each  $\pi = (\pi_0, \pi_1, \dots) \in \Pi$ ,  $\text{Prof}(\Delta_t \in D_1 | H_t) = \pi_t(D_1 | H_t)$  and  $\text{Prof}(X_{t+1} \in D_2 | H_{t-1}, \Delta_{t-1}, X_t = x, \Delta_t = a) = Q(D_2 | x, a)$ , for every  $D_1 \in B_A$  and  $D_2 \in B_S$ , and  $t=0, 1, 2, \dots$ .

For any Borel set  $X$ ,  $P(X)$  is the set of all probability measures on  $X$ . Then, for each  $\pi \in \Pi$  and initial state distribution  $\nu \in P(S)$ ,  $P_\pi^\nu$  is probability measure on  $\Omega$ , which can be defined in an obvious way, and  $E_\pi^\nu$  is the expectation with respect to  $P_\pi^\nu$ .

We define measurable functions on  $\Omega$  as follows:

$$(2.3) \quad \begin{aligned} \tilde{R}_T &:= \frac{1}{T} \sum_{t=0}^{T-1} r(X_t, \Delta_t) \quad (T \geq 1), \\ \tilde{R} &:= \liminf_{T \rightarrow \infty} \tilde{R}_T. \end{aligned}$$

$$(2.4) \quad \begin{aligned} \tilde{C}_T^i &:= \frac{1}{T} \sum_{t=0}^{T-1} c_i(X_t, \Delta_t) \quad (T \geq 1), \\ \tilde{C}^i &:= \limsup_{T \rightarrow \infty} \tilde{C}_T^i \quad (i = 1, \dots, k). \end{aligned}$$

For any  $(\nu, \pi) \in P(S) \times \Pi$ ,

$$(2.5) \quad R(\nu, \pi) := \text{ess} \cdot \inf \tilde{R} (= \sup \{a | P_\nu^\pi(\tilde{R} \geq a) = 1\})$$

$$(2.6) \quad C_i(\nu, \pi) := \text{ess} \cdot \sup \tilde{C}^i (= \inf \{a | P_\nu^\pi(\tilde{C}^i \leq a) = 1\}) \quad (i = 1, \dots, k),$$

where  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_k)$  is given.

Let

$$U_\alpha := \{(\nu, \pi) \in P(S) \times \Pi | C^i(\nu, \pi) \leq \alpha_i, i = 1, \dots, k\}$$

$$U_\alpha^{\text{RS}} := \{(\nu, \Phi^\infty) \in U_\alpha | (\nu, \Phi^\infty) \in P(S) \times \Pi_{\text{RS}}\}$$

In this report we mainly consider the following problem.

$$\begin{aligned} \text{Problem(A)} : \quad & \text{Maximum } R(\nu, \pi) \\ & \text{subject to } (\nu, \pi) \in U_\alpha \end{aligned}$$

$(\nu^*, \pi^*) \in U_\alpha$  will be called a constrained optimal pair if  $R(\nu^*, \pi^*) \geq R(\nu, \pi)$  for all  $(\nu, \pi) \in U_\alpha$ .

For any  $\epsilon > 0$ ,

$(\nu^*, \pi^*) \in U_\alpha$  is called a constrained  $\epsilon$ -optimal pair if  $R(\nu^*, \pi^*) \geq R(\nu, \pi) - \epsilon$  for all  $(\nu, \pi) \in U_\alpha$ .

In Section 3 we shall prove that a constrained optimal pair exists in  $U_\alpha^{\text{RS}}$  and in Section 4 give characterization of a constrained optimal pair.

### 3. Existence of optimal pair and related linear programmings

In this section, we transform Problem (A) given in the preceding section to LP equivalently and prove the existence of optimal pair by using compactness.

Let  $\{x_i\}$  be dense in  $S$  and define  $g_{ij} \in C(S)$  for  $i, j=1, 2, \dots$ , by

$$(3.1) \quad g_{ij}(x) = 2(1 - jd(x, x_i)) \vee 0,$$

where  $d$  is the metric defined in  $S$  and  $x \vee y = \max\{x, y\}$ .

Let

$$G = \{g_{ij} : i, j = 1, 2, \dots\}.$$

Then  $G$  is separating, i.e, whenever  $P_1, P_2 \in P(S)$  and

$$\int g dP_1 = \int g dP_2$$

for all  $g \in G$ , we have  $P_1 = P_2$  ([7]).

For  $\mu \in P(S \times A)$ ,  $h \in C(S \times A)$  we denote the integral as follows:

$$(3.2) \quad (h, \mu) := \int h(x, a) \mu(d(x, a)).$$

We can verify the following lemma by referring to the proof of Lemma 2.1 in ([8]).

**Lemma 3.1 .** For any positive  $\epsilon$ , and any  $(\nu, \pi) \in U_\alpha$ , there is  $\mu \in P(S \times A)$  such that:

- (i).  $(r, \mu) \geq R(\nu, \pi) - \epsilon$
- (ii).  $(c_i, \mu) \leq \alpha_i + \epsilon \quad (i = 1, 2, \dots, k)$
- (iii).  $\int g(x) \mu(d(x, a)) = \int \mu(d(x, a)) \int g(x') Q(dx'|x, a) \quad \text{for all } g \in G.$

Brief proof.

By the definition of  $R(\nu, \pi)$ ,  $C^i(\nu, \pi)$  and stability theorem [9] we can get the following:

$$(3.3) \quad P_\pi^\nu(\tilde{R} \geq R(\nu, \pi) - \epsilon) = 1$$

$$(3.4) \quad P_\pi^\nu(\tilde{C}_i \leq \alpha_i + \epsilon) = 1 \quad (i = 1, \dots, k)$$

$$(3.5) \quad \lim_{T \rightarrow \infty} \frac{\sum_{t=0}^{T-1} \{g(X_t) - E[g(X_t)|B_{t-1}]\}}{T} = 0 \quad \text{for all } g \in G, P_\pi^\nu - \text{almost surely}$$

From (3.3) to (3.5) there exists a sample path  $\omega \in \Omega$  such that  $\tilde{R}(\omega) \geq R(\nu, \pi) - \epsilon$ ,  $\tilde{C}_i(\omega) \leq \alpha_i + \epsilon$  ( $i = 1, \dots, k$ ) and (3.5) hold.

For this  $\omega \in \Omega$  we discuss the following empirical probability measure  $\mu_T \in P(S \times A)$ .

$$\mu_T(D) = \frac{\sum_{t=0}^{T-1} I_D(X_t, \Delta_t)}{T} \quad (T \geq 0), \quad \text{for all } D \in B_{S \times A},$$

where  $I_D$  is the indicator function of  $D$ .

From weak compactness of  $P(S \times A)$  there exists a sequence  $\{\mu_{T_j}\}$  such that (3.7) to (3.10) hold.

$$(3.7) \quad \lim_{j \rightarrow \infty} \left[ \int g(x) \mu_{T_j}(d(x, a)) - \int \mu_{T_{j-1}}(d(x, a)) \int g(x') Q(dx'|x, a) \right] = 0 \quad \text{for all } g \in G$$

$$(3.8) \quad \lim_{j \rightarrow \infty} (r, \mu_{T_j}) \geq R(\nu, \pi) - \epsilon$$

$$(3.9) \quad \lim_{j \rightarrow \infty} (c_i, \mu_{T_j}) \leq \alpha_i + \epsilon \quad (i = 1, 2, \dots, k)$$

$$(3.10) \quad \mu_{T_j} \rightarrow \mu \in P(S \times A), \text{ in the weak topology}$$

It is obvious that  $\mu$  in (3.10) satisfies (i),(ii),(iii) of lemma 3.1.

We now consider infinite linear programming as follows:

$$\begin{aligned} \text{LP*} \quad & \text{Maximum } (r, \mu) \\ & \text{subject to } \begin{aligned} & \text{(i)} \quad (c_i, \mu) \leq \alpha_i \quad (i = 1, \dots, k) \\ & \text{(ii)} \quad \int g(x) \mu(d(x, a)) = \int \mu(d(x, a)) \int g(x') Q(dx'|x, a) \\ & \text{(iii)} \quad \mu \in P(S \times A) \end{aligned} \end{aligned}$$

$P_F(S \times A)$  is the set of all  $\mu$  which satisfy the  $LP^*$  conditions (i),(ii),(iii).

In order to prove that Problem (A) is equivalent to  $LP^*$ , we need the following assumption 1, which remains operative thereafter.

Assumption 1.

For any  $\Phi^\infty \in \Pi_{RS}$ , the Markov chain induced by  $Q(\cdot|x, \Phi)$  satisfies the Doeblin condition and is one-ergodic, where  $Q(\cdot|x, \Phi) = \int Q(\cdot|x, a)\Phi(da|x)$ .

Under assumption 1 and lemma 3.1 we can verify the following theorem 3.1 by using weak compactness of  $P(S \times A)$  and ergodic theorem ([6]).

Theorem 3.1. For any  $(\nu, \pi) \in U_\alpha$ , there exists a  $\mu \in P_F(S \times A)$  such that for the decomposition  $\mu = \nu_s \times \Phi$ ,  $\nu_s \in P(S)$ ,  $\Phi \in P(A|S)$ , it holds:

- (i).  $(\nu_s, \Phi^\infty) \in U_\alpha$ ;
- (ii).  $R(\nu_s, \Phi^\infty) \geq R(\nu, \pi)$ ;
- (iii).  $R(\nu_s, \Phi^\infty) = (r, \mu)$ .

From the above result, the following follows.

Corollary 3.1. Problem (A) is equivalent to  $LP^*$

Since it is shown by compactness that  $LP^*$  has optimal solution (see [1]), the following holds from corollary 3.1.

Corollary 3.2. Optimal pair exists in  $U_\alpha^{RS}$ .

#### 4. State-wise mixed stationary policies

From corollary 3.1 we can get optimal pair or  $\epsilon$ -optimal pair of problem (A) by solving  $LP^*$ . In this section we give characterization of solutions of  $LP^*$ .

Let  $\mu_A \in P(A)$ . If there exists an integer  $l(l \geq 1)$ ,  $a_i \in A(i = 1, \dots, l)$  and  $p_i(i = 1, \dots, l)$  with  $p_i \geq 0$ ,  $\sum_{i=1}^l p_i = 1$  such that  $\mu_A(\{a_i\}) = p_i(i = 1, \dots, l)$ , we denote  $\mu_A$  by

$$(4.1) \quad \mu_A = \begin{pmatrix} a_1, a_2, \dots, a_l \\ p_1, p_2, \dots, p_l \end{pmatrix}.$$

Let  $\Phi \in P(A|S)$ , If there exists an integer  $l(l \geq 1)$ ,  $f_i \in B(S \rightarrow A)$  and  $p_i \in B(S \rightarrow [0, 1])$  ( $i = 1, \dots, l$ ) with  $p_i(x) \geq 0$ ,  $\sum_{i=1}^l p_i(x) = 1$  for all  $x \in S$ . such that

$$(4.2) \quad \Phi(\cdot|x) = \begin{pmatrix} f_1(x), f_2(x), \dots, f_l(x) \\ p_1(x), p_2(x), \dots, p_l(x) \end{pmatrix},$$

$\Phi$  is called an  $l$ -state-wise mixed kernel ( $l$ -s.m.k) and corresponding policy  $\Phi^\infty$  is called  $l$ -state-wise mixed stationary policy. Moreover, if  $\Phi$  is an  $l$ -s.m.k for some  $l \geq 1$ , we say  $\Phi$  is a s.m.k.

For  $l \geq 1$ , let

$$F^l := \{\mu \in P(S \times A) | \mu = \nu_s \times \Phi, \nu_s \in P(S), \Phi : l\text{-s.m.k}\}$$

$$F := \bigcup_{l=1}^{\infty} F^l$$

We note  $F^1$  represents the set of all non-randomized stationary policies.

Now, we need the following assumption 2, which remains operative thereafter.

Assumption 2.

The set of inner points of  $U_\alpha$  is non-empty.

Theorem 4.1. For any  $\epsilon > 0$ , an  $\epsilon$ -optimal solution of  $LP^*$  exists in  $F$ .

The proof of theorem 4.1 is given in a sequence of lemmas, some of interest.

Lemma 4.1.  $F$  is convex and if  $S$  is a finite set  $F^l$  is compact for every  $l$  ( $l \geq 1$ ).

The proof theorem 4.1 is done by induction on  $l$  which is the number of constraints in  $LP^*$ . For that, the following definitions are given.

$$P_F^{(m)}(S \times A) := \{\mu \in P(S \times A) | (c_i, \mu) \leq \alpha_i (i = 1, \dots, m) \text{ and } \mu \text{ satisfies (ii) of } LP^*\}$$

$$P_F^{(0)}(S \times A) := \{\mu \in P(S \times A) | \mu \text{ satisfies (ii) of } LP^*\}$$

We note  $P_F^{(m)}(S \times A)$  is compact and convex.

For each  $m$  ( $0 \leq m \leq k$ ), consider the following LP problem:

$$\begin{array}{ll} LP^{(m)} & \text{Maximum } (r, \mu) \\ & \text{subject to } \mu \in P_F^{(m)}(S \times A) \end{array}$$

In case that  $m = 0$ , the following holds (see [8]).

Lemma 4.2.  $LP^{(0)}$  has optimal solution in  $F^1$ .

Lemma 4.2 shows that theorem 3.1 is true for  $LP^{(0)}$ . Now for  $LP^{(m)}$  ( $1 \leq m \leq k$ ) it is supposed that theorem 4.1 is true. Here we use Lagrangian multiplier techniques [4].

Let

$$(4.3) \quad r_\lambda := r - \lambda c_{m+1} \quad (\lambda \geq 0)$$

We consider the following LP.

$$\begin{array}{ll} \text{LP}^{**} & \text{Maximum } (r_\lambda, \mu) \\ & \text{subject to } \mu \in P_F^{(m)}(S \times A) \end{array}$$

For any  $\epsilon > 0$ , it follows from the assumption of induction that an  $\epsilon$ -optimal solution of  $\text{LP}^{**}$  exist in  $F$ . Consequently, an optimal solution  $\mu^\lambda$  of  $\text{LP}^{**}$  exists in  $\bar{F}$  (the closure of  $F$ ).

Let

$$\begin{aligned} J^\lambda &:= (r_\lambda, \mu^\lambda) = (r, \mu^\lambda) - \lambda(c_{(m+1)}, \mu^\lambda); \\ R^\lambda &:= (\nu, \mu^\lambda); \\ K^\lambda &:= (c_{m+1}, \mu^\lambda). \end{aligned}$$

The following two lemmas can be verified by the a way as these in [4].

Lemma 4.3.  $J^\lambda, R^\lambda, K^\lambda$  are non-decrease functions of  $\lambda$ .

Lemma 4.4.  $\gamma := \inf\{\lambda : K^\lambda \leq \alpha_{m+1}\} < \infty$ .

Under the above preliminaries we know that if  $K^\gamma = \alpha_{m+1}$ ,  $\mu^\gamma$  is the optimal solution of  $\text{LP}^{(m+1)}$  by usually Lagarange multiplier way and if  $K^\gamma \neq \alpha_{m+1}$  Theorem 4.1 is shown to be true for  $\text{LP}^{(m+1)}$  by using lemma 4.1.

When  $S$  is a finite set, the following results are obtained by compact property of  $F^l$ .

Corollary 4.1. If  $S$  is finite set an optimal solution of  $\text{LP}^*$  exist in  $F^{2k}$ .

From corollary 3.1 the following fact can be easily obtained.

Corollary 4.2. For any  $\epsilon > 0$  there are  $\nu \in P(S)$  and s.m.k  $\Phi$  such that  $(\nu, \Phi^\infty)$  is an  $\epsilon$ -optimal pair. If  $S$  is finite set, there are  $\nu \in P(S)$  and 2k-s.m.k  $\Phi$  such that  $(\nu, \Phi^\infty)$  is an optimal pair.



# Bibliography

- [1] Anderson E.J. and Nash P. *Linear Programming in Infinite Dimensional Space*, Wiley, Chichester, 1987.
- [2] Arapostathis A. Borkar V.S. Fernandez-Gaucherand E. Ghosh M.K. and Marcus S.I. *Discrete time controlled Markov process with average cost criterion: a survey*, SIAM J. Control optim. vol.31, no.2, 282-344, 1993.
- [3] Bertsekas D.P. and Shreve S.E. *Stochastic optimal control the discrete times case*, Academic Press, New York, 1978.
- [4] Beutler F.J. and Ross K.W. *Optimal policies for controlled Markov chain with a constraint*, J. Math. Anal. Appl. 112(1985), 236-256.
- [5] Borkar V.S. *Topics in controlled Markov chain*, Pitman Research Notes in Math. No. 240, Longman Scientific and Technical, Harlow, 1991.
- [6] Doob J.L. *Stochastic Processes*, John Wiley, New York, 1953.
- [7] Ethier S.N. and Kurtz T.G. *Markov process characterization and convergence*, John Wiley, 1986.
- [8] Kurano M. *The existence of a minimum pair of state and policy for Markov decision processes under the hypothesis of Doeblin*, SIAM J. Control Optim. 27(1989), 296-307.
- [9] Loeve M. *Probability theory*, Second edition Van Nostrand, Princeton, NJ, 1960.
- [10] Ohnishi M. *Markov 決定過程に関する最近の話題*, Proceeding of the Fourth RAMP Symposium, 1992.
- [11] Ross K.W. *Randomized and path dependent policies for Markov decision processes with multiple constraints*, Open. Res. 37(1989) 474-477.
- [12] Ross K.W. and Varadarajan R. *Markov decision processes with sample path constraints: The communicating cases*, Open. Res. 37(1989), 780-790.
- [13] Ross K.W. and Varadarajan R. *Multiple Markov decision processes with a sample path constrain: A decomposition approach*, Math. O.R. 16(1991), 195-207.